# GENERAL EXPONENTIAL MODELS FOR DISCRETE OBSERVATIONS

## BY

## STEFFEN L. LAURITZEN

TECHNICAL REPORT NO. 17
MAY, 1974

DEPARTMENT OF STATISTICS
STANFORD  UNIVERSITY
STANFORD, CALIFORNIA

GENERAL EXPONENTIAL MODELS FOR

DISCRETE OBSERVATIONS


by

STEFFEN L. LAURITZEN

Institute of Mathematical Statistics
University of Copenhagen

and

Stanford University

TECHNICAL REPORT NO. 17

MAY, 1974

DEPARTMENT OF STATISTICS

STANFORD  UNIVERSITY

STANFORD, CALIFORNIA

## 1. Introduction and Summary

The purpose of the present paper is to illustrate the concept of an extreme family as defined by Lauritzen (1974) and to define a class of statistical models for discrete observations generalizing classical exponential families.

In the classical formulation, a discrete exponential family is a family of probability measures $(P_\theta, \theta \in \Theta)$, where the parameter space $\Theta$ is a subset of k-dimensional Euclidean vector space, the probability function being given by

$$P_\theta(x) = a(\theta)\, b(x)\, e^{\sum_{i=1}^{k} \theta_i t_i(x)} \tag{1.1}$$

where $x \in E$, a discrete set, $t_i$ are real valued functions and $\theta = (\theta_1, \ldots, \theta_k)$. If one observes independent identically distributed random variables $X_1, \ldots, X_n$ with the common probability for $X_i$ given by (1.1), the joint probability will be given by

$$P_\theta^{(n)}(x_1, \ldots, x_n) = a(\theta)^n \left( \prod_{j=1}^{n} b(x_j) \right) \cdot e^{\sum_{i=1}^{k} \theta_i \left( \sum_{j=1}^{n} t_i(x_j) \right)} \tag{1.2}$$

Now, $P_\theta^{(n)}$ is again an exponential family with the same parameter space as before. Somehow this is not a coincidence. If we try to look closer at the elements of the exponential family, we might understand this fact.

The function $b$ is a common reference measure defining the support of the measures $(P_\theta, \theta \in \Theta)$, and $a(\theta)$ is a normalizing constant. The functions

$(t_1,\ldots,t_k)$ are the sufficient statistics, and as an experiment is repeated, the sufficient statistics for the combined experiment is obtained as the <u>sum</u> of the sufficient statistics for the experiments in the repetition:

$$t_i^{(n)}(x_1,\ldots,x_n) = t_i(x_1) + \ldots + t_i(x_n) \ . \qquad (1.3)$$

The reason for this is that the function

$$g_\theta: (t_1,\ldots,t_k) \to e^{\sum_{i=1}^{k} \theta_i t_i}$$

is a homomorphism of the range space of $(t_1,\ldots,t_k)$ into the group $((0,\infty),\cdot)$:

$$g_\theta\left((t_1,\ldots,t_k) + (s_1,\ldots,s_k)\right) = g_\theta\left((t_1,\ldots,t_k)\right) \cdot g_\theta\left((s_1,\ldots,s_k)\right) \ . \qquad (1.4)$$

The idea in this paper is that most results about exponential families essentially are based on the above properties only. We shall therefore try to define a class of families of distributions via these properties.

If we again look at (1.1), (1.2) and (1.3) we see that we never substract. In fact, we only use that the algebraic operation + is associative and commutative. A set with a composition which is associative, commutative and has a unit is called a commutative monoid, cf. Bourbaki (1970). Commutative monoids are so simple that they are not studied very much. Therefore, we have to establish some of the simple results about these ourselves, which is done in Section 2.

Let us consider another family of distributions $(P_\theta, \theta \in \Theta)$, where $\Theta = \{1,2,\ldots,\}$ and

$$P_\theta\ (x) = \frac{1}{\theta} \cdot \chi_{\{1,\ldots,\theta\}}\ (x)\ ,\qquad\qquad (1.4)$$

where $x \in E = \{1,2,\ldots\}$ and $\chi_A$ is the indicator function of the set A, i.e.

$$\chi_A(x) = \begin{cases} 1 & \text{if } x \in A \\ 0 & \text{otherwise .} \end{cases}\qquad\qquad (1.5)$$

Consider $X_1,\ldots,X_n$ independent indentically distributed as above. Their joint probability is given by

$$P_\theta^{(n)}\ (x_1,\ldots,x_n) = \frac{1}{\theta^n}\chi_{\{1,\ldots,\theta\}}\ (\max\{x_1,\ldots,x_n\})\qquad\qquad (1.6)$$

The same situation as before is actually present if we replace $t(x_1)+ t(x_2)$ by $\max\{t\ (x_1),\ t\ (x_2)\}$. Just $g_\theta\ (x) = \chi_{\{1,\ldots,\theta\}}(x)$ can turn out to be zero this time whereas exponentials are always strictly positive. The support of the measures $(P_\theta, \theta \in \Theta)$ in this last example is varying with $\theta \in \Theta$, which is not the case in the first example.

In section 7, part I of Barndorff-Nielsen (1973), there is a detailed discussion of problems connected to maximum likelihood estimation in exponential families. The maximum likelihood estimator in regular canonical exponential families is shown to exist iff the observation happens to be so, that the value of the sufficient statistic falls within the interior of the convex hull of the support of the measures in the family, transformed by the sufficient statistics. This means that if the boundary of this convex hull has positive probability, one might very well get an observation from which

it is impossible to estimate.  To solve this problem it is proposed there
to make a suitable extension of the model, the extension being defined for
families where the set of possible values of the set of sufficient statistics
is assumed to be finite.  The extension is called the completion of an exponen-
tial family.

The measures in the <u>completion</u> of an exponential family have certainly
their support varying with the parameter, and the "fixed support" property
does therefore not seem to be essential to the nice results existing for
exponential families.

The families defined in the present paper are shown to be "complete" in
the sense that the maximum likelihood estimator of the parameters always
exist.

In section 3, the families are defined and some examples are discussed.
In section 4 we show the existence and uniqueness of the maximum likelihood
estimate of the unknown parameter in such families.

In section 5 we show that the family of Markov chains made up by sequences
of sufficient statistics from successive independent repetitions of an experi-
ment giving rise to a general exponential model, is in fact an extreme family
of Markov chains as defined by Lauritzen (1974).

In section 6 we shall briefly discuss the relation between the models
defined in the present paper and the completion of an exponential family as
defined by Barndorff-Nielsen (1973).


## 2.  Commutative Monoids

The algebraic structure of commutative monoids will play an essential role
in the present paper.  We shall quote the definition, cf. Bourbaki (1970).

<u>Definition 2.1</u>  Let  M  be a set and * a  composition rule on  M.  (M,*)
is said to be a <u>commutative monoid</u> if  *  is associative, commutative and
has a unit, i.e., if

$$\text{i)} \quad \forall a,b,c \in M: \quad a * (b*c) = (a*b) * c \quad,$$

$$\text{ii)} \quad \forall a,b \in M: \quad a * b = b * a \quad,$$

$$\text{iii)} \quad \exists a \in M: \; \forall a \in M: \quad e * a = a * e = a \quad.$$

As we shall only consider <u>commutative</u> monoids throughout this paper, we shall
just write "monoid" instead of "commutative monoid".  Examples of commutative
monoids are

1)  $(\underset{\sim}{N},+)$, where  $\underset{\sim}{N} = \{0,1,2,\ldots\}$.  Here  0  is the unit.

2)  $(\underset{\sim}{N},\vee)$, where  $x \vee y = \max\{x,y\}$.  The unit is  0.

3)  $(\underset{\sim}{N}\cup\{\infty\},\wedge)$, where  $x \wedge y = \min\{x,y\}$.  The unit is  $\infty$.

4)  $(\underset{\sim}{R}_+, \cdot)$, where  $\underset{\sim}{R}_+$  is the set of nonnegative real numbers.
    The unit is 1.

   Now, let  (M,*)  be a monoid.  Consider the set  $\hat{M}$  consisting of all
homomorphisms  $\xi:(M,*) \to (\underset{\sim}{R}_+,\cdot)$,  i.e., satisfying for all  a,b∈M

$$\xi(a) \; \xi(b) = \xi(a*b), \quad \xi(e) = 1, \tag{2.1}$$

where  e  is the unit in  (M,*).  If  $\xi_1, \xi_2 \in \hat{M}$,  the mapping  $\xi_1 \cdot \xi_2$  defined
by

$$\xi_1 \cdot \xi_2(a) = \xi_1(a) \; \xi_2(a) \tag{2.2}$$

is obviously in  $\hat{M}$  and it is a trivial exercise to verify that  $(\hat{M}, \cdot)$  is a
monoid with the unit being  $\xi_e$,  defined as

$$\xi_e(a) = 1 \quad \text{for all} \quad a \in M , \tag{2.3}$$

$(\hat{M}, \cdot)$ shall be called the <u>dual</u> monoid to $(M, *)$.

If we have two monoids $(M, *)$ and $(N, \circ)$ we can form the <u>product</u> of these

$$(M, *) \times (N, \circ) = (M \times N, \circledast), \tag{2.4}$$

where

$$(m_1, n_1) \circledast (m_2, n_2) = (m_1 * m_2, n_1 * n_2) . \tag{2.5}$$

This is again a monoid, the unit being $(e_M, e_N)$, where $e_M$ and $e_N$ are the units in $M$ and $N$ respectively. The dual to a product can easily be obtained from the duals to the elements in the product:

<u>Proposition 2.1</u>: <u>The homomorphisms of</u> $(M \times N, \circledast)$ <u>into</u> $(R_+, \cdot)$ <u>are exactly those of the form</u>

$$\xi(m, n) = \xi_M(m) \, \xi_N(n) ,$$

<u>where</u> $\xi_M \in \hat{M}$ <u>and</u> $\xi_N \in \hat{N}$ .

<u>Proof</u>: The equation

$$\xi(m_1 * m_2, \; n_1 \circ n_2) = \xi(m_1, n_1) \, \xi(m_2, n_2) \tag{2.6}$$

gives for $m_1 = m$, $m_2 = e_M$, $n_2 = n$, $n_1 = e_N$

$$\xi(m, n) = \xi(m, e_N) \cdot \xi(e_M, n) . \tag{2.7}$$

Now (2.6) for $n_1 = n_2 = e_N$ gives

$$\xi(m_1 * m_2, \; e_N) = \xi(m_1, \; e_N) \, \xi(m_2, \; e_N) \tag{2.8}$$

which means that $\xi(\cdot, e_N)$ must be in $\hat{M}$. Analogously one gets that $\xi(e_M, \cdot)$ must be in $\hat{N}$. So all homomorphisms must be of the form

$$\xi(m, n) = \xi_M(m) \ \xi_N(n) \ . \tag{2.9}$$

It is easy to see that functions of this form are in $\overline{\hat{M \times N}}$ . End of proof.

As mentioned in the introduction, the support of the measures in the families we consider may very often vary with the parameter. This will of course not be in a completely arbitrary fashion but in a fashion compatible with the algebraic structure of the sufficient statistics. To investigate this aspect, the following concept will be of relevance:

Definition 2.2 $F \subseteq M$ is said to be a face of $M$ if

   i) $F$ is a submonoid of $M$ and

   ii) $c \epsilon F \wedge c = a * b \implies a \epsilon F \wedge b \epsilon F$

The faces of $M$ are exactly the possible positivity regions for elements in $\hat{M}$:

Proposition 2.2: Let $F \subseteq M$ $F$ is a face of $M$ iff there is a $\xi \epsilon \hat{M}$ so that

$$F = \{a \epsilon M: \xi(a) > 0\} \ .$$

Proof: If $F = \{a \epsilon M: \xi(a) > 0\}$ for some $\xi \epsilon \hat{M}$, then

$$a \epsilon F \wedge b \epsilon F \implies \xi(a * b) = \xi(a) \ \xi(b) > 0 \ , \tag{2.10}$$

and as $\xi(e) = 1$, $F$ is a submonoid of $M$. If $c \epsilon F$ and $c = a * b$, then

$$0 < \xi(a * b) = \xi(a) \ \xi(b) \tag{2.11}$$

and hence $\xi(a)$ and $\xi(b)$ both must be positive, i.e., $a \epsilon F$ and $b \epsilon F$.

If on the other hand $F$ is a face of $M$, we can define

$$\xi(a) = \begin{cases} 1 & \text{if } a \in F \\ \\ 0 & \text{otherwise .} \end{cases} \qquad (2.12)$$

$\xi$ is easily seen to be a homomorphism and the result is proved. We also have

**Proposition 2.3:** $M$ is a face of $M$

The proof is obvious.

**Proposition 2.4.** If $(F_i)_{i \in I}$ is a family of faces of $M$, then $F = \bigcap_{i \in I} F_i$ is a face of $M$.

**Proof:** Immediate from the definition.

**Remark:** From propositions 2.3 and 2.4 it follows that for any $a \in M$ there is a unique smallest face of $M$, $F(a)$ so that $a \in F(a)$.

Propositions 2.1 and 2.2 enables us to establish a result about faces of product monoids:

**Proposition 2.5:** $F$ is a face of $M \times N$ iff $F = F_M \times F_N$, where $F_M$ and $F_N$ are faces of respectively $M$ and $N$.

**Proof:** According to proposition 2.1, all homomorphisms $\xi$ of $(M \times N, \circledast)$ into $(R_+, \cdot)$ are of the form

$$\xi(m,n) = \xi_M(n) \, \xi_N(n) \quad , \qquad (2.13)$$

where $\xi_M \in \hat{M}$ and $\xi_N \in \hat{N}$. Now

$$\left\{(m,n): \quad \xi(m,n) > 0\right\}$$

$$= \left\{(m,n): \quad \xi_M(m) > 0 \wedge \xi_N(n) > 0\right\}$$

$$= \left\{m: \xi_M(m) > 0\right\} \times \left\{n: \xi_N(n) > 0\right\} . \qquad (2.14)$$

The proposition 2.2 and equation (2.14) together yields the result.

Example 2.1: Let us consider the monoid $(\underset{\sim}{N},+)$. A homomorphism must satisfy $\xi(0) = 1$. Let $\xi(1) = \theta$, some non-negative real number. We must have

$$\xi(n) = \xi(1)^n = \theta^n \quad . \qquad (2.15)$$

It follows that the only faces of $(\underset{\sim}{N},+)$ are $\{0\}$ and $\underset{\sim}{N}$.

Example 2.2 Let us consider $(\underset{\sim}{N}, v)$. Let $n\epsilon\underset{\sim}{N}$ be fixed. The smallest face containing $n$ must contain all integers less than or equal to $n$ as

$$n \vee x = n \quad \text{if} \quad x \leq n \quad . \qquad (2.16)$$

On the other hand, $\{0,\ldots,n\}$ is obviously a face of $(\underset{\sim}{N}, v)$. Hence all faces of $(\underset{\sim}{N}, v)$ are $\underset{\sim}{N}$ itself and subsets of the form $\{0,\ldots,n\}$ for some $n\epsilon\underset{\sim}{N}$. Now let $\xi \in \underset{\sim}{\hat{N}}$ be positive exactly on $\{0,\ldots,n\}$; it must satisfy

$$\xi(x \vee n) = \xi(x) \ \xi(n) = \xi(n) \quad \text{for} \quad x \leq n. \qquad (2.17)$$

As $\xi(n)$ is strictly positive, we get

$$\xi(x) = 1 \quad \text{for } x \leq n \quad , \qquad (2.18)$$

and hence that the <u>only</u> homomorphisms of $(\mathbb{N},+)$ are indicator functions of faces,

$$\xi(x) = \chi_{\{0,\ldots,n\}}(x) \qquad\qquad (2.19)$$

for some $n \in \mathbb{N} \cup \{\infty\}$.

<u>Example 2.3</u>: If we now form the product $(\mathbb{N},\vee) \times (\mathbb{N},+)$ it follows from proposition 2.1 that all homomorphisms are of the form

$$\xi(x,y) = \chi_{\{0,\ldots,n\}}(x) \cdot \theta^y \qquad\qquad (2.20)$$

for some $n \in \mathbb{N} \cup \{\infty\}$ and $\theta \geq 0$. From proposition 2.5, we get that the faces of this monoid are

$$\{0,\ldots,n\} \times \mathbb{N}, \quad n \in \mathbb{N}$$

$$\{0,\ldots,n\} \times \{0\}, \quad n \in \mathbb{N}$$

$$\mathbb{N} \times \mathbb{N} \quad .$$

## 3. General Exponential Models

In the following we shall consider an at most denumerable set $E$, a monoid $(M,*)$ and a mapping $t: E \to M$. We shall think of $E$ as the sample space and of $t$ as a sufficient statistic. Let $M_1 = t(E)$ and define recursively

$$M_n = M_1 * M_{n-1}, \quad \text{for } n = 2,3,\ldots \quad . \qquad\qquad (3.1)$$

This is done for the following reason: if we make $n$ independent observations of a random variable on $E$, we shall assume that the sufficient statistic will be

$$t^{(n)}(x_1,\ldots,x_n) = t(x_1) * \cdots * t(x_n) \ , \tag{3.2}$$

and hence $M_n = t^{(n)}(E^n)$ .

We shall assume that we can infer the size of the experiment from the statistic, or, in other words, that

$$M_m \bigcap M_n = \emptyset, \quad \text{whenever} \quad n \neq m \ . \tag{3.3}$$

For convenience we don't want $M$ to be bigger than necessary, hence we assume that

$$M = \bigcup_{n=0}^{\infty} M_n \ , \tag{3.4}$$

where $M_o$ satisfies

$$M_0 * M_n = M_n \quad \text{for all} \quad n \in \underset{\sim}{N} \ . \tag{3.5}$$

Let $\nu$ be a $\sigma$-finite measure on $E$ so that $\nu(x)$ is positive for all $x \in E$. Let $\hat{M}_\nu$ denote the <u>normalized dual</u> to $M$:

$$\hat{M}_\nu = \left\{ \xi \in \hat{M} : \sum_{x \in E} \nu(x)\, \xi(t(x)) = 1 \right\} \tag{3.6}$$

and assume that $\hat{M}_\nu$ is non-empty.

A statistical model for a random variable $X$ taking values in $E$ is a family $\mathcal{P}$ of probability measures on $E$.

<u>Definition 3.1</u>: $\mathcal{P}$ is said to be a <u>general exponential model</u> if there exists $M, t$ and $\nu$ as above, so that

$$P \epsilon \mathcal{P} \iff \exists \, \xi \in \hat{M}_\nu : \quad P\{X = x\} = \nu(x) \, \xi(t(x)) \quad .$$

Let us first see in what sense this looks like a "classical" exponential model. Suppose we have observed $n$ independent random variables from the above distribution. The joint probability function is

$$P_\xi^{(n)} (x_1, \ldots, x_n) = \left( \prod_{i=1}^n \nu(x_i) \right) \cdot \xi(t(x_1)) \cdots \xi(t(x_n))$$

$$= \left( \prod_{i=1}^n \nu(x_i) \right) \xi(t(x_1) * \cdots * t(x_n)) \qquad (3.7)$$

as $\xi \epsilon \hat{M}$. If we compare (3.7) with (1.2) in the introduction, we note, that $\nu$ plays the same role as $b$, the common reference measure. The statistic $t$ corresponds to $(t_1, \ldots, t_k, n)$, i.e. the sufficient statistic plus a "counting variable" indicating the size of the experiment. $\xi \epsilon \hat{M}_\nu$ corresponds to the function

$$(t_1, \ldots, t_k, n) \to a(\theta)^n \, e^{\sum_{i=1}^k \theta_i t_i} \qquad (3.8)$$

so the normalizing constant $a(\theta)$ is taken into $\xi$ and the experiment size into the statistic $t$.

The above defined models differ from the exponential models in several respects. First the range space of the statistic is a monoid instead of a subset of a vector space, the parameter space is the normalized dual of this monoid instead of a subset of a vector space, and there is no assumption of anything like finite dimension. Furthermore we shall see that in general the support of the measures in the family will depend on the parameter $\xi$,

as these will not always be positive. As derived in the previous sections, the possible positivity regions for $\xi$ will be the faces of the monoid $(M, *)$. From (3.7) and the Neyman factorization theorem it immediately follows that

$$t^{(n)}(X_1, \ldots, X_n) = t(X_1) * \cdots * t(X_n) \tag{3.9}$$

is sufficient for the parameter $\xi$ from observation of $X_1, \ldots, X_n$.

The relation between these models and the classical exponential models should hopefully be more apparent from the examples below.

Example 3.1 (The Bernoulli Distribution)

Let $E = \{0, 1\}$ and $\nu(0) = \nu(1) = 1$. Let

$$(M, *) = (\underset{\sim}{N}, +) \times (\underset{\sim}{N}, +) \ . \tag{3.10}$$

We have

$$M = \bigcup_{n=0}^{\infty} M_n, \quad \text{where} \quad M_n = \{(x, y): x + y = n\} \ . \tag{3.11}$$

Let $t(1) = (1, 0)$ and $t(0) = (0, 1)$. The elements of $\hat{M}$ are all of the form

$$F_{\theta, \eta}(x, y) = \theta^x \eta^y, \quad \theta \geq 0, \quad \eta \geq 0. \tag{3.12}$$

We immediately get that

$$F_{\theta, \eta} \epsilon \hat{M}_\nu \iff \theta^1 \eta^0 + \theta^0 \eta^1 = 1$$

$$\iff \eta = 1 - \theta \ . \tag{3.13}$$

Hence, the model

$$P_\theta \{X = x\} = \nu(x) \, F_{\theta, 1-\theta}(t(x)) = \begin{cases} \theta & \text{if } x = 1 \\ \\ 1-\theta & \text{if } x = 0 \ , \end{cases} \tag{3.14}$$

where $0 \leq \theta \leq 1$, is a general exponential model. The difference between this model and the classical exponential family version of the Bernoulli distribution is that $\theta = 0$ and $\theta = 1$ are included in the model.

Example 3.2 (The Poisson distribution).

Let $E = \underset{\sim}{N}$ and $\nu(x) = \frac{1}{x!}$ . Let

$$(M, *) = (\underset{\sim}{N}, +) \times (\underset{\sim}{N}, +). \tag{3.15}$$

We have

$$M = \bigcup_{n=0}^{\infty} M_n, \quad \text{where} \quad M_n = \{(x,y); \ y = n\}. \tag{3.16}$$

Let $t(x) = (x,1)$. We get

$$F_{\theta,\eta} \in \hat{M}_\nu \iff \sum_{x=0}^{\infty} \frac{1}{x!} \theta^x \eta = 1$$

$$\iff \eta = e^{-\theta} \tag{3.17}$$

Hence, the model

$$P_\theta\{X = x\} = \nu(x) \, F_{\theta, e^{-\theta}}(t(x)) = \frac{\theta^x}{x!} e^{-\theta} \ , \tag{3.18}$$

where $\theta \geq 0$ is a general exponential model. Again the inclusion of $\theta = 0$ is the only difference between this and the classical approach.

So far, the examples considered have basically been exponential models in the classical sense apart from adding some degenerate distributions. The following examples show that the models in fact can be quite different from the classical exponential models.

Example 3.3 (The uniform distribution).

Let $E = \{1,2,\ldots\}$ and $\nu(x) = 1$ for all $x \in E$. Let

$$(M,*) = (E,\nu) \times (\underset{\sim}{N},+) . \tag{3.19}$$

We have

$$M = \bigcup_{n=0}^{\infty} M_n, \quad \text{where} \quad M_n = \{(x,y): y = n\} \tag{3.20}$$

Let $t(x) = (x,1)$. The elements of $\hat{M}$ are all of the form

$$F_{\theta,\eta} (x,y) = \chi_{\{1,\ldots,\theta\}} (x)\eta^y, \quad \theta \in E, \quad \eta \geqq 0 . \tag{3.21}$$

We have

$$F_{\theta,\eta} \in \hat{M}_\nu \Longleftrightarrow \sum_{x=1}^{\infty} \eta \chi_{\{1,\ldots,\theta\}}(x) = 1$$

$$\Longleftrightarrow \eta = \frac{1}{\theta} . \tag{3.22}$$

Hence, the model

$$P_\theta \{X = x\} = \nu(x) F_{\theta,\frac{1}{\theta}} (t(x)) = \frac{1}{\theta} \cdot \chi_{\{1,\ldots,\theta\}} (x) , \tag{3.23}$$

where $\theta = 1,2,\ldots,$ is a general exponential model.

Combining examples 3.1 - 3.3 we take the following:

<u>Example 3.4</u>  (Doubly truncated geometric distribution with unknown truncation points).  Let $E = \underset{\sim}{N}$  and  $\nu(x) = 1$  for all  $x \epsilon E$.  As our monoid  $(M, *)$  we choose the submonoid of

$$(M', *) = (\underset{\sim}{N}, \nu) \times (\underset{\sim}{N} \cup \{\infty\}, \wedge) \times (\underset{\sim}{N}, +) \times (\underset{\sim}{N}, +) \quad , \tag{3.24}$$

given by  $M = \bigcup_{n=0}^{\infty} M_n,$  where

$$M_0 = \{(0, \infty, 0, 0)\}, \quad M_1 = \{(x, x, x, 1): x \epsilon \underset{\sim}{N}\} \quad , \tag{3.25}$$

and  $M_n$  is recursively defined as

$$M_n = M_1 * M_{n-1} \quad \text{for} \quad n = 2, 3, \dots \tag{3.26}$$

Let  $t(x) = (x, x, x, 1)$.  The elements of  $\hat{M}$  are all of the form

$$F_{\theta, \eta, \lambda, \mu} (x, y, z, n) =$$

$$\chi_{\{0, \dots, \theta\}}(x) \ \chi_{\{\eta, \dots, \infty\}}(y) \ \lambda^z \ \mu^n \quad , \tag{3.27}$$

where  $\theta \epsilon \underset{\sim}{N}$,  $\eta \epsilon \underset{\sim}{N} \cup \{\infty\}$,  $\lambda \geq 0$  and  $\mu \geq 0$ .

We get

$$F_{\theta, \eta, \lambda, \mu} \epsilon \hat{M}_\nu \iff \sum_{x=\eta}^{\theta} \lambda^x \mu = 1$$

$$\iff \eta \leq \theta, \ \lambda = 1 \quad \text{and} \quad \frac{1}{\mu} = \theta - \eta + 1$$

$$\text{or} \quad \eta \leq \theta, \ \lambda \neq 1 \quad \text{and} \quad \frac{1}{\mu} = \frac{\lambda^{\theta+1} - \lambda^\eta}{\lambda - 1} \quad . \tag{3.28}$$

Hence the model

$$P_{\theta,\eta,\lambda} \{X = x\} = \phi(\lambda)\, \lambda^x \chi_{\{\eta,\ldots,\theta\}} (x) \quad , \tag{3.29}$$

where $\lambda \geqq 0$, $\theta \geqq \eta$, $\theta$, $\eta \in \underset{\sim}{N}$ and

$$\phi(\lambda) = \begin{cases} \dfrac{1}{\theta-\eta+1} & \text{if } \lambda = 1 \\[2ex] \dfrac{\lambda-1}{\lambda^{\theta+1}-\lambda^{\eta}} & \text{if } \lambda \neq 1 \quad , \end{cases} \tag{3.30}$$

is a general exponential model.

Finally we shall consider an example, where the general exponential model is different from a classical one in the sense of infinite-dimensionality of the parameter space.

Example 3.5 (The completely free distribution). Let $E$ be any denumerable set, and $\nu(x) = 1$ for all $x \in E$. Let $(M,*) = \underset{\sim}{{}^E N}$, consisting of all mappings $f$ from $E$ to $\underset{\sim}{N}$ with finite support, i.e., where $\{x: f(x) \neq 0\}$ is finite. The composition rule is pointwise addition

$$(f*g)\,(x) = f(x) + g(x) \quad . \tag{3.31}$$

We have the partitioning of $M = \bigcup_{n=0}^{\infty} M_n$, where

$$M_n = \left\{ f \in {}^E\underset{\sim}{N}: \sum_{x \in E} f(x) = n \right\} \quad . \tag{3.32}$$

If we let $t(x) = \chi_{\{x\}}$, we can see, that the sufficient reduction of a sample of size $n$ becomes the "frequency table", i.e., $t^{(n)}(x_1,\ldots,x_n)$ is the function in ${}^E\underset{\sim}{N}$, having the value $n_x$ in $x$ iff $x$ occurs exactly $n_x$ times in the sample $(x_1,\ldots,x_n)$. $\hat{M}$ consists of the elements

$$g_\theta \ (f) = \prod_{x \in E} \theta(x)^{f(x)} \quad , \tag{3.33}$$

where $\theta$ is any mapping from $E$ into the non-negative real numbers. We have

$$g_\theta \in \hat{M}_\nu \iff \sum_{x \in E} \left( \prod_{y \in E} \theta(y)^{\chi_{\{x\}}(y)} \right) = 1$$

$$\iff \sum_{x \in E} \theta(x) = 1 \quad . \tag{3.34}$$

Hence, the model

$$P_\theta \ \{X = x\} = \nu(x) \ g_\theta \ (\chi_{\{x\}}) = \theta(x) \quad , \tag{3.35}$$

where $\theta$ satisfies

$$\left. \begin{array}{l} \theta(x) \geq 0 \quad \text{for all} \quad x \in E \quad \text{and} \\[2mm] \displaystyle\sum_{x \in E} \theta(x) = 1 \end{array} \right\} \quad , \tag{3.36}$$

is a general exponential model. Other examples could be generated ad libitum.

## 4. Estimation in general exponential models

We shall consider the following estimation problem:

Let $X_1, \ldots, X_n$ be independent and identically distributed on $E$ with

$$P_\xi \ \{X = x\} = \nu(x) \ \xi(t(x)) \quad , \tag{4.1}$$

where $\nu$ and $t$ are known and as in the previous section and $\xi \in \hat{M}_\nu$ is unknown. Our sample space is $E^n$, the parameter space is $\hat{M}_\nu$ and the likelihood function becomes

$$L(x_1, \ldots, x_n, \ \xi) = \prod_{x=1}^{n} \nu(x_i) \quad \xi(t(x_1) * \cdots * t(x_n)). \tag{4.2}$$

As mentioned earlier, $t^{(n)}$ given by

$$t^{(n)}(x_1,\ldots,x_n) = t(x_1)*\cdots*t(x_n) \tag{4.3}$$

is sufficient for $\xi$ and $\hat{\xi}$ is clearly a maximum likelihood estimator of $\xi$ iff

$$\hat{\xi}(t_o) = \sup_{\xi \in \hat{M}_\nu} \xi(t_o) \quad , \tag{4.4}$$

where $t_o = t(x_1)*\cdots*t(x_n)$ .

In the following we shall establish the existence and uniqueness of $\hat{\xi}$ for any $n$ and $x_1,\ldots,x_n$.

First we prove a lemma:

<u>Lemma 4.1</u>   Let $\hat{M}_\nu^* = \left\{ \xi \in \hat{M}: \sum_{x \in E} \nu(x)\, \xi(t(x)) \leqq 1 \right\}$ . $\hat{M}_\nu^*$ <u>is compact in the pointwise topology</u>.

<u>Proof:</u>   Let $\xi_1, \xi_2, \ldots$ be a sequence of elements in $\hat{M}_\nu^*$. As $[0, \infty]^M$ is compact, we can always find a subsequence $\xi_{n_1}, \xi_{n_2}, \ldots$ so that for any $s \in M$,

$$\xi_{n_i}(s) \xrightarrow[i \to \infty]{} \xi(s) \quad , \tag{4.5}$$

where $0 \leqq \xi(s) \leqq \infty$.

We have to show that this limit $\xi$ in fact is an element of $\hat{M}_\nu^*$. From Fatou's lemma, we get

$$\sum_{x \in E} \nu(x)\, \xi(t(x)) \leqq \liminf_{i \to \infty} \sum_{x \in E} \nu(x)\, \xi_{n_i}(t(x)) \leqq 1 \tag{4.6}$$

We shall now just prove that $\xi(s) < +\infty$ and that

$$\xi(s * t) = \xi(s)\, \xi(t) \; . \tag{4.7}$$

But as $t(E) = M_1$ and $\nu(x) > 0$ for all $x \in E$, (4.6) gives that $\xi(s) < +\infty$ for all $s \in M_1$. Now, if $s \in M_n = M_1 * \cdots * M_1$, where $M_1$ appears $n$ times, we have

$$\xi(s) = \lim_{i \to \infty} \xi_{n_i}(s) = \lim_{i \to \infty} (\xi_{n_i}(s_1) \cdots \xi_{n_i}(s_n)) = \xi(s_1) \cdots \xi(s_n) \quad , \tag{4.8}$$

where $s_1, \ldots, s_n \in M$, and $s_1 * \cdots * s_n = s$. This gives that $\xi(s) < \infty$ for all $s \in M$ since

$$\xi(e) = \lim_{i \to \infty} \xi_{n_i}(e) = 1, \tag{4.9}$$

and also that $\xi \in \hat{M}$. The lemma is proved.

We can now show the existence of the maximum likelihood estimate for any $n, x_1, \ldots, x_n$:

Proposition 4.1: For all $s \in M$, there is a $\hat{\xi} \in \hat{M}_\nu$, so that $\hat{\xi}(s) = \sup\limits_{\xi \in \hat{M}_\nu} \xi(s)$.

Proof:

As $\hat{M}_\nu^*$ is compact and the mapping $\xi \to \xi(s)$ is continuous, there is a $\xi^* \in \hat{M}_\nu^*$ so that

$$\xi^*(s) = \sup_{\xi \in \hat{M}_\nu^*} \xi(s) \; . \tag{4.10}$$

But if

$$\sum_{x \in E} \nu(x)\, \xi^*(t(x)) = c < 1 \tag{4.11}$$

then $\hat{\xi}: M \to \underset{\sim}{R}_+$ defined as

$$\hat{\xi}(s) = \xi^*(s) \left(\frac{1}{c}\right)^n \quad \text{for} \quad s \in M_n \tag{4.12}$$

is in $\hat{M}_\nu^*$ and $\hat{\xi}(s) > \xi^*(s)$, which is a contradiction. Hence we must have $c = 1$, $\hat{\xi}(s) = \xi^*(s)$, $\hat{\xi} \in \hat{M}_\nu$ and

$$\hat{\xi}(s) = \sup_{\xi \in \hat{M}_\nu} \xi(s) \quad , \tag{4.13}$$

which was to be proved.

Next we prove the uniqueness of the maximum likelihood estimate.

<u>Proposition 4.2</u>: If $\hat{\xi}_1(s_0) = \hat{\xi}_2(s_0) = \sup_{\xi \in \hat{M}_\nu} \xi(s_0)$ , then $\hat{\xi}_1 = \hat{\xi}_2$.

<u>Proof</u>: For $s \in M$ let

$$\xi(s) = \sqrt{\xi_1(s) \, \xi_2(s)} \quad . \tag{4.14}$$

Define $\hat{\xi}: M \to \underset{\sim}{R}_+$ by

$$\hat{\xi}(s) = \frac{\xi(s)}{\left(\sum_{x \in E} \nu(x) \, \xi(t(x))\right)^k} \quad \text{for} \quad s \in M_k \, . \tag{4.15}$$

Obviously $\hat{\xi} \in \hat{M}_\nu$. If $\hat{\xi}_1 = \hat{\xi}_2$ for all $s \in M_1$, $\hat{\xi}_1 = \hat{\xi}_2$ for all $s \in M$. Cauchy-Schwarz inequality gives

$$\sum_{x \in E} \nu(x) \, \xi(t(x)) \leqq \left(\sum_{x \in E} \nu(x) \, \hat{\xi}_1(t(x))\right) \left(\sum_{x \in E} \nu(x) \, \hat{\xi}_2(t(x))\right) = 1, \tag{4.16}$$

as $\hat{\xi}_1$ and $\hat{\xi}_2$ are in $\hat{M}_\nu$. We therefore have

$$\hat{\xi}_1 \neq \hat{\xi}_2 \implies \sum_{x \in E} \nu(x) \, \xi(t(x)) < 1. \tag{4.17}$$

But then

$$\hat{\xi}(s_0) > \sqrt{\hat{\xi}_1(s_0) \, \hat{\xi}_2(s_0)} = \hat{\xi}_1(s_0) = \hat{\xi}_2(s_0) \quad, \tag{4.18}$$

which is a contradiction. Hence $\hat{\xi}_1 = \hat{\xi}_2 = \hat{\xi}$ which was to be proved.

The next result giving some more detailed information about the maximum likelihood estimate should be compared to the results in section 7, part I of Barndorff-Nielsen (1973).

<u>Proposition 4.3</u>: <u>The positivity region of</u> $\hat{\xi}$ <u>where</u> $\hat{\xi}(s_0) = \sup_{\xi \in \hat{M}_\nu} \xi(s_0)$ <u>is</u> <u>exactly the face</u> $F(s_0)$.

<u>Proof</u>: As $\hat{\xi}(s_0) > 0$, we have from proposition 2.2 that

$$M^+(\hat{\xi}) = \{s \in M : \hat{\xi}(s) > 0\} \supseteq F(s_0) . \tag{4.19}$$

Suppose that $M^+(\hat{\xi}) \neq F(s_0)$, i.e., there is an $s'$ in $M^+(\hat{\xi}) \backslash F(s_0)$. Then $s' \in M_n$ for some $n$ and $s' = s_1 * s_2 * \cdots * s_n$. As $M^+(\hat{\xi})$ is a face, we then have

$$s_1, \ldots, s_n \in M^+(\hat{\xi}) \cap M_1 . \tag{4.20}$$

At least one of them, say $s_1$, must be outside $F(s_0)$ since $F(s_0)$ is a submonoid and the sum is outside $F(s_0)$. Now let

$$\xi'(s) = \begin{cases} \hat{\xi}(s) & \text{for } s \in F(s_0) \\ 0 & \text{otherwise} \end{cases} \tag{4.21}$$

and define $\xi: M \to \underset{\sim}{R}_+$ by

$$\xi(s) = \frac{\xi'(s)}{\left(\sum_{x \in E} \nu(x)\xi'(t(x))\right)^n} \quad \text{for} \quad s \in M_n. \tag{4.22}$$

Clearly $\xi \in \hat{M}_\nu$ and

$$\sum_{x \in E} \nu(x) \ \xi'(t(x)) < 1. \tag{4.23}$$

Therefore $\xi(s_0) > \hat{\xi}(s_0)$, which is a contradiction. Hence we must have

$$M^+(\hat{\xi}) = F(s_0) \ , \tag{4.24}$$

which was to be proved.

So, the support of the estimated measure is closely tied to the way the observations can occur. If $s_0$ is observed after $n$ experiments, $t(x_i)$ <u>must</u> be in the smallest face containing $s_0$ for all $i = 1,\ldots,n$ as

$$s_0 = t(x_1)_* \cdots {}_*t(x_n) \ . \tag{4.25}$$

The estimate contains this information as the support of $P_{\hat{\xi}}$ is reduced to the subset of $E$, where $t(x) \in F(s_0)$.

## 5. Extreme Families of Random Walks on Monoids

First we introduce the definition of an extreme family of Markov chains as given by Lauritzen (1974).

Let $(E_n, n=1,2,\ldots)$ be a family of discrete, at most denumerable spaces and $Q = (Q_{mn})_{m \leq n}$ a family of matrices with elements $q_{mn}(x,y)$, $x \in E_m$, $y \in E_n$, satisfying

$$\left. \begin{array}{l} q_{mn}(x,y) \geq 0, \qquad \displaystyle\sum_{x \in E_m} q_{mn}(x,y) = 1 \\[4mm] \text{and} \\[2mm] Q_{mn} Q_{np} = Q_{mp} \qquad \text{for} \quad m \leq n \leq p \ . \end{array} \right\} \tag{5.1}$$

$\mathfrak{M}(Q)$ denotes the set of sequences of probability measures $\mu = (\mu_n, \; n = 1,2,\ldots)$ so that $\mu_n$ is a probability measure on $E_n$ and

$$\mu_m = Q_{mn} \, \mu_n \quad \text{for all} \quad m \leqq n. \tag{5.2}$$

$\mathfrak{M}(Q)$ is a convex set and $\mathfrak{E}(Q)$ shall mean the extreme points of $\mathfrak{M}(Q)$. The family of Markov chains on $\Pi_{n=1}^{\infty} E_n$ defined by the initial distributions

$$P^\mu \{X_1 = x\} = \mu_1(x) \tag{5.3}$$

and the transition probabilities for $m \leq n$

$$P^\mu \{X_n = y | X_m = x\} = \begin{cases} q_{mn}(x,y) \, \dfrac{\mu_n(y)}{\mu_m(x)} & \text{for} \quad \mu_m(x) \neq 0 \\[2mm] \mu_n(y) & \text{otherwise} \end{cases} \tag{5.4}$$

where $\mu$ takes all values in $\mathfrak{E}(Q)$, is called the extreme family generated by $Q$.

For all $\mu \in \mathfrak{M}(Q)$, the matrices $Q_{mn}$ define the "backward conditional probabilities", i.e.

$$P^\mu \{X_m = x | X_n = y\} = q_{mn}(x,y) \quad \text{for} \quad m \leqq n \tag{5.5}$$

and $\mu_n$ the marginal distributions of $X_n$, i.e.

$$P^\mu \{X_n = x\} = \mu_n(x) \; . \tag{5.6}$$

Now consider the sequence of spaces $M_n$, $n=1,2,\ldots$ where $M = U_{n=0}^{\infty} M_n$ corresponding to a general exponential model. Let

$$Y_n = t(X_1) * \cdots * t(X_n), \tag{5.7}$$

where $X_1,\ldots,X_n$ are independent and identically distributed as

$$P_\xi \{X = x\} = \nu(x)\, \xi(t(x)) \; , \tag{5.8}$$

where $\xi \in \hat{M}_\nu$ is unknown. Let

$$\alpha(s) = \sum_{x \in E: t(x)=s} \nu(x) \; . \tag{5.9}$$

We have $\alpha(s) > 0$ for all $s \in M_1$. Define the $n$'th convolution $\alpha^{*n}$ of $\alpha$ as $\alpha^{*1}(s) = \alpha(s)$ and

$$\alpha^{*n}(s) = \sum_{a*b=s} \alpha(a)\, \alpha^{*(n-1)}(b) \quad \text{for} \quad m=2,3,\ldots \tag{5.10}$$

We have $\alpha^{*n}(s) > 0$ for all $s \in M_n$.

$Y_1,\ Y_2,\ldots$ forms a Markov chain on $\Pi_{n=1}^{\infty} M_n$ and we have for $m \leqq n$ and $P_\xi \{Y_n=y\} > 0$

$$P_\xi \{Y_m = x \mid Y_n = y\} = \frac{P_\xi\{Y_n=y \mid Y_m=x\} \cdot P_\xi\{Y_m=x\}}{P_\xi\{Y_n=y\}}$$

$$= \frac{\left(\displaystyle\sum_{a:\, a*x=y} \alpha^{*(n-m)}(a)\, \xi(a)\right) \alpha^{*m}(x)\, \xi(x)}{\alpha^{*n}(y)\, \xi(y)}$$

$$= \frac{\alpha^{*m}(x)}{\alpha^{*n}(y)} \sum_{a:\, a*x=y} \alpha^{*(n-m)}(a) \; . \tag{5.11}$$

We shall now consider the system of backward conditional distributions $(Q_{mn})_{m \leqq n} = Q$ with elements $q_{mn}(x,y)$ $x \in M_m, y \in M_n$ given by

$$q_{mn}(x,y) = \frac{\alpha^{*m}(x)}{\alpha^{*n}(y)} \sum_{a:\, a*x=y} \alpha^{*(n-m)}(a) \tag{5.12}$$

We shall find $\xi(Q)$ and in fact show that

$$\mathcal{E}(Q) = \left\{ \mu: \; \exists \, \xi \in \hat{M}_\nu: \; \mu_n(x) = \alpha^{*n}(x) \, \xi(x) \right\} \qquad (5.13)$$

i.e. exactly the family of Markov chains made up by sequences of sufficient statistics from successive repetitions of experiments giving rise to random variables following a general exponential model.

First we need a lemma. For $\mu = (\mu_n, n=1,2,\ldots) \in \mathcal{M}(Q)$, $x \in M_k$, $k=1,2,\ldots$ define the sequence $T_{x,k} \, \mu$ by

$$T_{x,k} \, \mu_n(a) = \begin{cases} \dfrac{\mu_{n+k}(a*x)}{\mu_k(x)} \; \dfrac{\alpha^{*k}(a) \, \alpha^{*n}(x)}{\alpha^{*(n+k)}(a*x)} & \text{if } \mu_k(x) > 0 \\[3mm] \mu_n(a) & \text{otherwise} \end{cases} \qquad (5.14)$$

<u>Lemma 4.1</u>  $\quad \mu \in \mathcal{M}(Q) \Rightarrow T_{x,k} \mu \in \mathcal{M}(Q)$ .

<u>Proof:</u>

Clearly, if $\mu_k(x) = 0$, $T_{x,k} \, \mu = \mu$ and hence $T_{x,k} \, \mu \in \mathcal{M}(Q)$.

If $\mu_k(x) \neq 0$, we have

$$\sum_{a \in M_n} T_{x,k} \, \mu_n(a) = \frac{1}{\mu_k(x)} \cdot \sum_{b \in M_{n+k}} \sum_{a : a*x=b} \frac{\alpha^{*k}(x) \, \alpha^{*n}(a)}{\alpha^{*(n+k)}(b)} \, \mu_{n+k}(b)$$

$$= \frac{1}{\mu_k(x)} \cdot \sum_{b \in M_{n+k}} q_{n,n+k}(x,b) \, \mu_{n+k}(b) = 1 \qquad (5.15)$$

as $\mu$ was known to be in $\mathcal{M}(Q)$.

Further, we get

$$\sum_{b \in M_n} q_{mn}(a,b) \, T_{x,k} \, \mu_n(b)$$

$$= \sum_{b \in M_n} \sum_{c : c*a=b} \frac{\alpha^{*m}(a) \, \alpha^{*(n-m)}(c)}{\alpha^{*n}(b)} \; \frac{\mu_{n+k}(b*x)}{\mu_k(x)} \; \frac{\alpha^{*n}(b) \, \alpha^{*k}(x)}{\alpha^{*(n+k)}(b*x)}$$

$$= \frac{\alpha^{*m}(a) \; \alpha^{*k}(x)}{\alpha^{*(m+k)}(a*x)} \cdot \frac{1}{\mu_k(x)} \cdot \sum_{b \in M_n} \sum_{c:c*a=b} \frac{\alpha^{*(m+k)}(a*x) \; \alpha^{*(n-m)}(c)}{\alpha^{*(n+k)}(b*x)} \mu_{n+k}(b*x)$$

$$(5.16)$$

Now

$$\{c:c*a = b\} \subseteq \{c:c*a*x = b*x\} \tag{5.17}$$

and

$$\{b*x: \; x \in M_n\} \subseteq M_{n+k}, \tag{5.18}$$

so we have the inequality

$$\sum_{b \in M_n} q_{mn}(a,b) \; T_{x,k} \; \mu_n(b)$$

$$\leq \frac{\alpha^{*m}(a) \; \alpha^{*k}(x)}{\alpha^{*(m+k)}(a*x)} \frac{1}{\mu_k(x)} \cdot \sum_{d \in M_{n+k}} \sum_{c:c*(a*x)=d} \frac{\alpha^{*(m+k)}(a*x) \; \alpha^{*(n-m)}(c)}{\alpha^{*(n+k)}(d)} \mu_{n+k}(d)$$

$$= \frac{\alpha^{*m}(a) \; \alpha^{*k}(x)}{\alpha^{*(m+k)}(a*x)} \frac{1}{\mu_k(x)} \sum_{d \in M_{n+k}} q_{m+k, \; n+k}(a*x,d) \; \mu_{n+k}(d)$$

$$= T_{x,k} \; \mu_m(a) \; , \tag{5.19}$$

or in short,

$$T_{x,k} \; \mu_m(a) \geq \sum_{b \in M_n} q_{mn}(a,b) \; T_{x,k} \; \mu_n(b) \; . \tag{5.20}$$

But by (5.15), both sides of (5.20) add up to one when summing over $a \in M_m$, and hence we must have equality and therefore $T_{x,k} \; \mu \in \mathfrak{M}(Q)$, which was to be proved.

## Proposition 4.1

$$\mu \in \mathcal{E}(Q) \iff \exists\, \xi \in \hat{M}_\nu : \quad \mu_n(x) = \alpha^{*n}(x)\,\xi(x) \ .$$

Or, in words, the extreme family generated by $Q$ consist of "random walks",

$$Y_n = t(X_1) * \cdots * t\,(X_n) \tag{5.21}$$

where the $t(X)$'s are independent and identically distributed with the distribution of $Y_1 = t(X_1)$ given by

$$\mu_1(x) = \alpha(x)\,\xi(x) \quad \text{for some} \quad \xi \in \hat{M}_\nu \ . \tag{5.22}$$

Proof: The proof consists of the following steps. First we use lemma 4.1 to obtain a representation of any $\mu \in \mathcal{M}(Q)$ as a convex combination of other elements $(T_{x,k}\,\mu)$ in $\mathcal{M}(Q)$. If $\mu$ then is extreme, $\mu$ must be equal to these other elements, which gives us an equation. This equation is essentially the homomorphism equation and we can then establish "$\Rightarrow$". To prove "$\Leftarrow$" we show that a proper mixture of homomorphisms, cannot be a homomorphism.

Suppose now that $\mu \in \mathcal{M}(Q)$ is extreme, i.e. $\mu \in \mathcal{E}(Q)$. We note that the equation for $\mu \in \mathcal{M}(Q)$.

$$\mu_n(a) = \sum_{b \in M_{n+k}} \sum_{c:c*a=b} \frac{\alpha^{*n}(a)\,\alpha^{*k}(c)}{\alpha^{*(n+k)}(a*c)}\,\mu_{n+k}\,(a*c) \tag{5.23}$$

implies that

$$\mu_n(a) = 0 \Rightarrow \mu_{n+k}\,(a*c) = 0 \ , \tag{5.24}$$

as $q_{n,n+k}(a,b) > 0$ for all $b \in M_{n+k}$. Hence (5.23) can be rearranged to

$$\mu_n(a) = \sum_{b \in M_{n+k}} \sum_{c:c*a=b} \mu_k(c) \cdot T_{c,k}\,\mu_n(a) \tag{5.25}$$

This gives $\mu$ as a convex combination of $T_{c,k}\,\mu$, $c \in M_k$ for all $k = 1,2,\ldots$ and as $\mu$ was supposed to be extreme,

$$\mu_n(a) = T_{c,k}\,\mu_n(a) \quad \text{for all} \quad k = 1,2,\ldots \quad \text{and} \quad c \in M_k. \quad (5.26)$$

Thus, for all $c \in M_k$, $k = 1,2,\ldots$ so that $\mu_k(c) > 0$, we must have

$$\frac{\mu_{n+k}(a*c)}{\alpha^{*(n+k)}(a*c)} = \frac{\mu_n(a)}{\alpha^{*n}(a)} \frac{\mu_k(c)}{\alpha^{*k}(c)} \qquad (5.27)$$

If we let

$$h_n(a) = \frac{\mu_n(a)}{\alpha^{*n}(a)} \quad , \qquad (5.28)$$

(5.27) becomes

$$h_{n+k}(a*c) = h_n(a)\,h_k(c) \qquad (5.29)$$

But as

$$\mu_k(c) = 0 \implies \mu_{n+k}(a*c) = 0 \qquad (5.30)$$

(5.29) must hold for <u>all</u> $n$, $k$, $a \in M_n$, $c \in M_k$. As $M_n \cap M_m = \emptyset$ for $m \neq n$, we can define a mapping $\xi$ from $M$ to $\underset{\sim}{R}_+$ by

$$\xi(a) = h_n(a) \quad \text{for} \quad a \in M_n \ , \qquad (5.31)$$

and by (5.29) $\xi \in \hat{M}$.

If $\mu$ is extreme, we then have

$$\frac{\mu_n(a)}{\alpha^{*n}(a)} = \xi(a) \iff \mu_n(a) = \alpha^{*n}(a)\,\xi(a). \qquad (5.32)$$

As $\mu$ is a probability, we have

$$\sum_a \mu_n(a) = \sum_a \alpha^{*n}(a)\ \xi(a) = 1 \quad, \tag{5.33}$$

i.e., that $\xi \in \hat{M}_\nu$. We have proved "$\Rightarrow$".

Now suppose that

$$\mu_n^{\xi_0}(x) = \alpha^{*n}(x)\ \xi_0(x) \quad \text{for some} \quad \xi_0 \in \hat{M}_\nu \ . \tag{5.34}$$

All elements in $\mathcal{M}(Q)$ are mixtures of the elements in $\mathcal{E}(Q)$. It follows from what we proved before that the set of sequences

$$\{\mu^\xi,\ \xi \in \hat{M}_\nu\} \ , \tag{5.35}$$

where

$$\mu_n^\xi(x) = \alpha^{*n}(x)\ \xi(x) \quad , \tag{5.36}$$

contains $\mathcal{E}(Q)$. A forteriori any $\mu \in \mathcal{M}(Q)$ can be represented as a mixture of elements of the form (5.36). This is in particular true for $\mu^{\xi_0}$. Hence, there is a probability measure $P$ on $\hat{M}_\nu$ so that for all $x \in M$

$$\alpha^{*n}(x)\ \xi_0(x) = \int_{\xi \in \hat{M}_\nu} \alpha^{*n}(x)\ \xi(x)\ d\ P(\xi) \ . \tag{5.37}$$

But (5.37) is equivalent to

$$\xi_0(x) = \int_{\xi \in \hat{M}_\nu} \xi(x)\ d\ P(\xi) \quad \text{for all} \quad x \in M \tag{5.38}$$

Using the homomorphism property, we have

$$\left(\int_{\xi \in \hat{M}_{\nu}} \xi(x) \, d \, P(\xi)\right)^2 = (\xi_o(x))^2 = \xi_o(x * x)$$

$$= \int_{\xi \in \hat{M}_{\nu}} \xi(x * x) \, d \, P(\xi) = \int_{\xi \in \hat{M}_{\nu}} (\xi(x))^2 \, d \, P(\xi) \quad (5.39)$$

But (5.39) implies that $P\{\xi_o\} = 1$ and hence that $\xi_o$ is extreme. The proof is complete.

## 6. Additional Comments

The families defined in the present paper are sometimes identical to the completion of a regular canonical exponential family as defined by Barndorff-Nielsen (1973).

Let $T$ be a finite subset of the k-dimensional integer lattice let $T_o = \{\underline{o}\}$ and define $T_n$ by

$$T_1 = T \quad \text{and} \quad T_n = T + T_{n-1} \quad \text{for} \quad n = 2,3,\ldots \quad (6.1)$$

Let $(M, *)$ be the monoid

$$M = \{(t,n): \ t \in T_n, \ n \in \underset{\sim}{N}\} \quad , \quad (6.2)$$

with the composition

$$(s, m) * (t, n) = (s + t, m + n) \quad . \quad (6.3)$$

If we have a general exponential model on a space $E$ with $t(x) = (g(x), 1)$, where $g$ is a mapping from $E$ onto $T$, this model can be identified with the completion of the canonical exponential family generated by $\nu$ and $g$ in exactly the same way as in Martin-Löf (1973). This situation is present in example 3.1 of the present paper.

If $g(E) = T$ contains more than integer lattice points, this is <u>not</u> necessarily the case as the following example shows.

<u>Example 6.1</u>   Let $E = \{(0,0), (1,0), (0,1), (\sqrt{2}/4, 1/2)\}$.  Let $M = \cup_{n=0}^{\infty} M_n$, where $M_0 = \{(0,0,0)\}$ and

$$M_n = \{(x,y,n): (x,y) \in E_n \text{ and } n \in \underset{\sim}{N}\} \tag{6.4}$$

where $E_n$ is recursively defined as

$$E_1 = E \text{ and } E_n = E + E_{n-1} \text{ for } n = 2,3,\ldots \tag{6.5}$$

and the composition on $(M, *)$ is defined as

$$(x,y,n) * (x', y', n') = (x + x', y + y', n + n'). \tag{6.6}$$

Let $\nu(x,y) = 1$ for all $(x,y)$ in $E$ and

$$t(x,y) = (x,y,1) \quad . \tag{6.7}$$

The completion of the exponential family generated by $\nu$ and $t$ would consist of all probability measures with support equal to $E$ and the probability degenerate in $(1,0)$, $(0,0)$ and $(0,1)$. Because $(\sqrt{2}/4, 1/2)$ is in the interior of the convex hull of $E$, no probability in the family would be degenerate at this point.

The subset $F$ of $M$ given by

$$F = \left\{\left(n \frac{\sqrt{2}}{4}, \frac{n}{2}\right) : n \in \underset{\sim}{N}\right\} \tag{6.8}$$

is obviously a face of $M$. Hence the general exponential model corresponding to $\nu$, $t$ and $M$ <u>will contain</u> the probability degenerate in $(\sqrt{2}/4, 1/2)$.

This example illustrates the essential difference between the completions defined by Barndorff-Nielsen (1973) and the models in this paper: the "completions" are defined via geometrical concepts in $\underset{\sim}{R}^k$ or via topological considerations whereas the general exponential models are derived via algebraic structure in the statistics, thus letting the actual observations play a more prominent role. If one after $n$ experiments in the above example obtains the value of the statistic to be $(n\sqrt{2}/4, n/2)$, this must be because $(\sqrt{2}/4, 1/2)$ was observed $n$ times. This is reflected in the estimated probability measure, which will be degenerate at $(\sqrt{2}/4, 1/2)$ as can be seen from proposition 4.3.

## Acknowledgments

I am grateful to all colleagues at the Institute of Mathematical
Statistics, University of Copenhagen for several inspiring discussions on
the topics of the present paper. My sincere thanks are due to Ann Mitchell,
Imperial College, Søren Johansen and Niels Keiding, Copenhagen, and
T. W. Anderson, Stanford, for reading various versions of this paper and
suggesting important alterations in the manuscript. Finally, this paper
would never exist if Hans Brøns, Copenhagen ever missed an opportunity to
point out, that algebraic properties of sufficient statistics are essential
features of the models they correspond to and vice versa.

## References

Barndorff-Nielsen, O. (1973). Exponential Families and Conditioning. Copenhagen, Wiley.

Bourbaki, N. (1970). Algèbre, Ch. I. Hermann, Paris.

Lauritzen, S. L. (1974). Sufficiency, Prediction and Extreme Models. Scandinavian Journal of Statistics, 1.

Martin-Löf, P. (1973). Repetitive Structures. Department of Mathematical Statistics. University of Stockholm. (Mimeographed report).

TECHNICAL REPORTS

OFFICE OF NAVAL RESEARCH CONTRACT N00014-67-A-0112-0030 (NR-042-034)

1. "Confidence Limits for the Expected Value of an Arbitrary Bounded Random Variable with a Continuous Distribution Function," T. W. Anderson, October 1, 1969.

2. "Efficient Estimation of Regression Coefficients in Time Series," T. W. Anderson, October 1, 1970.

3. "Determining the Appropriate Sample Size for Confidence Limits for a Proportion," T. W. Anderson and H. Burstein, October 15, 1970.

4. "Some General Results on Time-Ordered Classification," D. V. Hinkley, July 30, 1971.

5. "Tests for Randomness of Directions against Equatorial and Bimodal Alternatives," T. W. Anderson and M. A. Stephens, August 30, 1971.

6. "Estimation of Covariance Matrices with Linear Structure and Moving Average Processes of Finite Order," T. W. Anderson, October 29, 1971.

7. "The Stationarity of an Estimated Autoregressive Process," T. W. Anderson, November 15, 1971.

8. "On the Inverse of Some Covariance Matrices of Toeplitz Type," Raul Pedro Mentz, July 12, 1972.

9. "An Asymptotic Expansion of the Distribution of "Studentized" Classification Statistics," T. W. Anderson, September 10, 1972.

10. "Asymptotic Evaluation of the Probabilities of Misclassification by Linear Discriminant Functions," T. W. Anderson, September 28, 1972.

11. "Population Mixing Models and Clustering Algorithms," Stanley L. Sclove, February 1, 1973.

12. "Asymptotic Properties and Computation of Maximum Likelihood Estimates in the Mixed Model of the Analysis of Variance," John James Miller, November 21, 1973.

13. "Maximum Likelihood Estimation in the Birth-and-Death Process," Niels Keiding, November 28, 1973.

14. "Random Orthogonal Set Functions and Stochastic Models for the Gravity Potential of the Earth," Steffen L. Lauritzen, December 27, 1973.

15. "Maximum Likelihood Estimation of Parameters of an Autoregressive Process with Moving Average Residuals and Other Covariance Matrices with Linear Structure," T. W. Anderson, December, 1973.

16. "Note on a Case-Study in Box-Jenkins Seasonal Forecasting of Time Series," Steffen L. Lauritzen, April, 1974.

OFFICE OF NAVAL RESEARCH CONTRACT N00014-67-A-0112-0030 (NR-042-034)

17.  "General Exponential Models for Discrete Observations," Steffen L. Lauritzen, May, 1974.

| REPORT DOCUMENTATION PAGE | | READ INSTRUCTIONS BEFORE COMPLETING FORM |
|---|---|---|
| 1. REPORT NUMBER<br><br>Technical Report No. 17 | 2. GOVT ACCESSION NO. | 3. RECIPIENT'S CATALOG NUMBER |
| 4. TITLE (and Subtitle)<br><br>General Exponential Models for Discrete Observations | | 5. TYPE OF REPORT & PERIOD COVERED<br><br>Technical Report |
| | | 6. PERFORMING ORG. REPORT NUMBER |
| 7. AUTHOR(s)<br><br>Steffen L. Lauritzen | | 8. CONTRACT OR GRANT NUMBER(s)<br><br>N00014-67-A-0112-0030 |
| 9. PERFORMING ORGANIZATION NAME AND ADDRESS<br><br>Department of Statistics<br>Stanford University<br>Stanford, California 94305 | | 10. PROGRAM ELEMENT, PROJECT, TASK AREA & WORK UNIT NUMBERS<br><br>NR-042-034 |
| 11. CONTROLLING OFFICE NAME AND ADDRESS<br>Office of Naval Research<br>Statistics & Probability Program Code 436<br>Arlington, Virginia 22217 | | 12. REPORT DATE<br>April 1974 |
| | | 13. NUMBER OF PAGES<br>40 |
| 14. MONITORING AGENCY NAME & ADDRESS(If different from Controlling Office) | | 15. SECURITY CLASS. (of this report)<br><br>Unclassified |
| | | 15a. DECLASSIFICATION/DOWNGRADING SCHEDULE |

16. DISTRIBUTION STATEMENT (of this Report)

Reproduction in whole or in part is permitted for any purpose of the United States Government. Distribution is unlimited.

17. DISTRIBUTION STATEMENT (of the abstract entered in Block 20, if different from Report)

18. SUPPLEMENTARY NOTES

19. KEY WORDS (Continue on reverse side if necessary and identify by block number)

Commutative monoids, exponential families, extreme models, maximum likelihood estimation, sufficiency.

20. ABSTRACT (Continue on reverse side if necessary and identify by block number)

(see next page)

A class of models generalizing exponential families is defined via the algebraic structure of the sufficient statistics.  The maximum likelihood estimate for the unknown parameter is shown to exist and be unique.

The sequence of sufficient statistics from successive repetitions of experiments corresponding to a general exponential model is shown to form an extreme family of Markov chains as defined by Lauritzen (1974).